

# CALCOLO SCIENTIFICO

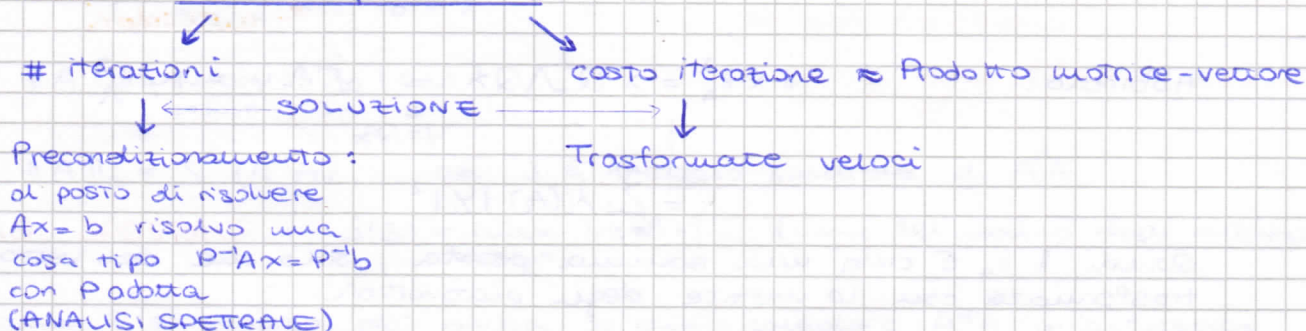
17/10/2012

①

## PROGRAMMA:

- Premesse e richiami
- SVD (decomposizione ai valori singolari)
- Metodi iterativi → Krylov (cg, gmres)  
→ Multigrid (geometrica e algebrica)

Riduzione costo computazionale:



## 1. PREMESSE e RICHIAMI

1.1 NORME MEMO: vogliamo risolvere

$$Ax=b$$

alla fine

def: **NORMA VETTORIALE**  $\|\cdot\|: \mathbb{C}^n \rightarrow \mathbb{R}^+ \cup \{0\}$  t.c.

- $\|x\| \geq 0$  e vale  $= 0$  sse  $x=0$
- $\|\alpha x\| = |\alpha| \|x\| \quad \forall \alpha \in \mathbb{C} \quad \forall x \in \mathbb{C}^n$
- $\|x+y\| \leq \|x\| + \|y\|$

**ESEMPLI:**  $\|x\|_2 = \left( \sum_{i=1}^n |x_i|^2 \right)^{1/2}$

$$\|x\|_\infty = \max_{i=1, \dots, n} |x_i|$$

$$\|x\|_1 = \left( \sum_{i=1}^n |x_i| \right)$$

**PROP<sub>1</sub>:**  $\|\cdot\|$  è una funzione uniformemente continua:

$$\forall \varepsilon > 0 \quad \exists \delta = \delta(\varepsilon) : \forall x, y \in \mathbb{C}^n \quad \|x-y\| \leq \delta \Rightarrow |\|x\| - \|y\|| \leq \varepsilon$$

Questo mi dice che se misuro vettori con numeri piccoli o numeri grandi non cambia la differenza delle norme.

Detto meglio, la norma della differenza non dipende dalla grandezza dei numeri.

**PROP<sub>2</sub>:** EQUIVALENZA TOPOLOGICA:  $\forall \|\cdot\|_\alpha, \|\cdot\|_\beta: \mathbb{C}^n \rightarrow \mathbb{R}^+ \cup \{0\}$  norme vettoriali  
 $\exists \alpha, \beta \in \mathbb{R} \quad 0 < \alpha \leq \beta$  t.c.

$$\alpha \|x\|_\beta \leq \|x\|_\alpha \leq \beta \|x\|_\beta \quad \forall x \in \mathbb{C}^n$$

Questo ci dice che alla fine possiamo usare la norma più comoda da calcolare tanto più o meno è uguale.

**PROP<sub>3</sub>:**  $\forall x \in \mathbb{C}^n \quad \|x\|_\infty \leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty$

$$\|x\|_2 \leq \|x\|_1 \leq \sqrt{n} \|x\|_2$$

$$\|x\|_\infty \leq \|x\|_1 \leq n \|x\|_\infty$$

def: **NORMA ENERGIA:**  $A$  sdp (simmetrica def pos / hermit.),  $A = G^H G$   
con  $G$  non singolare

Fatti Cholesky  
↑  
3 perche  
sdp

$$\|\cdot\|_{G^H G}: \mathbb{C}^n \rightarrow \mathbb{R}^+ \cup \{0\}, \quad \|x\|_{G^H G} := \|Gx\|_2 = (x^H G^H G x)^{1/2} = (x^H A x)^{1/2}$$

**OSS:** Siccome  $A = G^H G$  posso definirla dando  $A$  o dando  $G$  ed è lo stesso.



Se  $A = I$  allora  $\|x\|_A^2 = \|x\|_2^2 = \sum_{i=1}^n |x_i|^2$

Se  $A = D$  diagonale def pos ( $D = (d_i)$  con  $d_i > 0$ )  
 $\|x\|_A^2 = x^H D x = \sum_{i=1}^n d_i |x_i|^2$   
 (summa pesata)

Se  $A$  è generica sdf (l.h.f) posso sempre scrivere  
 $A = Q^H \Lambda Q$  con  $Q$  unitaria (matrice autovettori) i.e.  $Q Q^H = Q^H Q = I$   
 e  $\Lambda = \text{diag}(\lambda_i(A)) > 0$  (autovalori)

Abbiamo:  $\|x\|_A^2 = x^H Q^H \Lambda Q x = y^H \Lambda y = \|y\|_\Lambda^2$   
 $y = Qx$   
 $= \sum_{i=1}^n \lambda_i(A) |y_i|^2$

Quindi  $\|\cdot\|_A$  è cmq una summa pesata, solo che dei vettori trasformati con la matrice degli autovettori.

Usare una norma energia serve per dare importanza ad alcune componenti rispetto alle altre.

A ogni norma possiamo associare un prodotto scalare:

$\|\cdot\|_2 \rightsquigarrow \langle x, y \rangle = x^H y = \sum_{i=1}^n \bar{x}_i y_i$

$\|\cdot\|_{G^H G} \rightsquigarrow \langle x, y \rangle_{G^H G} = \langle x, G^H G y \rangle = \langle Gx, Gy \rangle$   
 $(\Rightarrow \|x\|_{G^H G}^2 = \langle Gx, Gx \rangle)$

def: **NORMA MATRICIALE**:  $\|\cdot\|: \mathbb{C}^{n \times n} \rightarrow \mathbb{R}^+ \cup \{0\} + \dots$

- $\|X\| \geq 0 \quad \forall X \in \mathbb{C}^{n \times n}$  e vale = 0 sse  $X = 0$
- $\|\alpha X\| = |\alpha| \|X\| \quad \forall \alpha \in \mathbb{C} \quad \forall X \in \mathbb{C}^{n \times n}$
- $\|X + Y\| \leq \|X\| + \|Y\|$
- $\|XY\| \leq \|X\| \|Y\|$  p. submoltiplicativa (non vale l'="a pensare dato che magari  $XY = 0$  ma  $X, Y \neq 0$ )

Come prima vogliamo le  $\text{PROP}_1$  e  $\text{PROP}_2$  dato che seguono dalle prime 3 proprietà.

def:  $\|\cdot\|_*$  è una **NORMA MATRICIALE COMPATIBILE** con il-norma vettoriale se:

$\|Ax\| \leq \|A\|_* \|x\| \quad \forall x \in \mathbb{C}^n, \forall A \in \mathbb{C}^{m \times n}$

def: **NORMA MATRICIALE INDOTTA** da una norma vettoriale:

$\|A\| = \max_{\|x\|=1} \|Ax\| = \sup_{x \in \mathbb{C}^n, \|x\|=1} \frac{\|Ax\|}{\|x\|}$

(si verifica che è effettivamente una norma)

**PROP<sub>4</sub>**: La norma indotta è una norma compatibile

dim: Se  $x = 0$   $Ax = 0 \Rightarrow 0 = 0$  ok

Se  $x \neq 0$  allora  $\|Ax\| = \|x\| \frac{\|Ax\|}{\|x\|} = \|x\| \left\| \frac{Ax}{\|x\|} \right\|$

poniamo  $y = \frac{x}{\|x\|} \rightarrow \|y\| = 1$

$\|Ax\| = \|x\| \|Ay\| \leq \|x\| \max_{\|y\|=1} \|Ay\| = \|x\| \|A\|$



**PROP<sub>5</sub>:** la n.m. indotta è la più piccola norma tra le norme compatibili

$$\text{dim: } \|A\| = \max_{\|x\|=1} \|Ax\| \leq \max_{\|x\|=1} \|A\|_1 \|x\| = \|A\|_1 \quad \text{e-cost}$$

**PROP<sub>6</sub>:** Indichiamo le norme indotte come le norme vettoriali, allora valgono:

$$\|A\|_1 = \max_{j=1, \dots, n} \left( \sum_{i=1}^n |a_{ij}| \right)$$

N.B. queste non sono definizioni, le def sarebbero:

$$\|A\|_\infty = \max_{i=1, \dots, n} \left( \sum_{j=1}^n |a_{ij}| \right)$$

$$\|A\|_1 = \max_{\|x\|_1=1} \|Ax\|_1 \quad \text{ecc...}$$

$$\|A\|_2 = \sqrt{\lambda}(A^H A) \quad \text{ove } \lambda \text{ è il raggio spettrale di } A^H A$$

**def: RAGGIO SPETTRALE:**  $\lambda(B) = \max_{i=1, \dots, n} |\lambda_i(B)|$  (max del modulo degli autovalori)

**OSS:** la norma 2 è proprio più brutta, perché calcolare  $A^H A$  in aritmetica floating point può essere un problema, e poi anche calcolare gli autovalori, quindi conviene usare per es. la  $\| \cdot \|_1$  che è molto più comoda e equivalente

**OSS:** se  $A = A^H \Rightarrow \|A\|_1 = \|A\|_\infty$

$$\|A\|_2 = \sqrt{\lambda}(A^2) = \sqrt{\lambda(A^2)} = \lambda(A)$$

perché se  $Ax = \lambda x$  allora  $A^2 x = \lambda^2 x$

se  $A$  è hdp  $\Rightarrow \|A\|_2 = \lambda_{\max}(A)$  (posso togliere il modulo tutti sono positivi)

**PROP:**  $\forall A \in \mathbb{C}^{n \times n}$

$$\frac{1}{\sqrt{n}} \|A\|_\infty \leq \|A\|_2 \leq \sqrt{n} \|A\|_\infty$$

$$\frac{1}{\sqrt{n}} \|A\|_1 \leq \|A\|_2 \leq \sqrt{n} \|A\|_1$$

$$\max_{i,j} |a_{ij}| \leq \|A\|_2 \leq n \max_{i,j} |a_{ij}|$$

**def: NORMA di FROBENIUS:**  $\|A\|_F := \left( \sum_{i,j=1}^n |a_{ij}|^2 \right)^{1/2} = [\text{tr}(A^H A)]^{1/2}$

**N.B.:** Non è una norma matriciale indotta!

$$\text{dim: } \|I\|_F = \sqrt{n} > 1 = \max_{\|x\|=1} \|Ix\|$$

Ogni norma indotta (si vede da qui) è t.c.  $\|I\| = 1$  e quindi  $\| \cdot \|_F$  non è indotta.

**PROP:**  $A \in \mathbb{C}^{m \times n}$ ,  $U \in \mathbb{C}^{n \times n}$  unitaria allora

$$\|UA\|_2 = \|AU\|_2 = \|UAU^H\|_2 = \|A\|_2$$

e idem per la norma  $\| \cdot \|_F$ .

$$\text{dim: } \|B\|_2^2 = \lambda(B^H B)$$

$$\|B\|_F^2 = \text{tr}(B^H B) = \sum_{i=1}^n \lambda_i(B^H B)$$

$$(UA)^H UA = A^H U^H UA = A^H A$$

$$(AU)^H (AU) = U^H A^H AU \sim_{\text{simile}} A^H A$$

$$(UAU^H)^H (UAU^H) = UA^H U^H UA U^H = UA^H AU^H \sim_{\text{simile}} A^H A$$

quindi hanno lo stesso raggio spettrale e la stessa traccia.

Conclusione: le matrici unitarie ci piacciono!



## 1.2 FORMA CANONICA di JORDAN

**TEO:**  $A \in \mathbb{C}^{n \times n}$  è diagonalizzabile  $A = S \Lambda S^{-1}$ ,  $AS = \Lambda S$  sse ha  $n$  autovettori linearmente indipendenti

Non diagonalizzabile se  $\mu_3 < \mu_2 \leftarrow$  mult. algebrica  
 mult. geometrica

**TEO:**  $A \in \mathbb{C}^{n \times n}$  con autovalori  $\lambda_i$   $i=1, \dots, n$  con molteplicità  $\mu_2(i), \mu_3(i)$

$A \sim J = \begin{bmatrix} J_1 & & \\ & \ddots & \\ & & J_p \end{bmatrix}$  diagonale a blocchi

con  $p = \#$  autovalori distinti, e  $J_i$  diagonale a blocchi

$J_i = \begin{bmatrix} C_{\lambda_i^{(i)}}^{(i)} & & \\ & \ddots & \\ & & C_{\mu_j(i)}^{(i)} \end{bmatrix} \in \mathbb{C}^{\mu_2(i) \times \mu_2(i)}$

con  $C_j^{(i)}$  **BLOCCO di JORDAN**

$C_j^{(i)} = \begin{bmatrix} \lambda_i & 1 & 0 \\ & \lambda_i & 1 \\ 0 & & \lambda_i \end{bmatrix} \in \mathbb{C}^{\nu_j(i) \times \nu_j(i)}$

con  $\sum_{j=1}^{\mu_2(i)} \nu_j(i) = \mu_2(i)$

## 1.3 FORMA CANONICA di SCHUR

**TEO:**  $A \in \mathbb{C}^{n \times n}$  con autovalori  $\lambda_1, \dots, \lambda_n$

$\Rightarrow \exists U, T \in \mathbb{C}^{n \times n}$   $U$  unitaria ( $UU^H = U^H U = I$ );  
 $T$  triang. sup + c.  $\text{diag}(T) = (\lambda_1, \dots, \lambda_n)$

t.c.  $A = UTU^H$

dim: Per induzione

- $n=1$   
Prendo  $U = [1]$  ed è ovvio
- $n > 1$   
Considero  $x_1$  autovettore normalizzato (in norma 2) risp. a  $\lambda_1$  autoval  
 $S := \langle x_1 \rangle$

Considero  $S^\perp$  con base ortonormale  $(y_2, \dots, y_n)$

Costruiamo  $Q = [x_1 | y_2 | \dots | y_n]$

Abbiamo che  $Q$  è unitaria e  $Q^H x_1 = e_1$

$Q^H Q = \begin{bmatrix} x_1^H x_1 & x_1^H y_2 & \dots & x_1^H y_n \\ x_2^H y_2 & y_2^H y_2 & & \\ \vdots & & \ddots & \\ x_n^H y_n & & & y_n^H y_n \end{bmatrix}$   
 1 perché  $x_1$  è normalizzato  
 0 perché  $y_2, \dots, y_n$  è base di  $S^\perp$

Poniamo  $B = Q^H A Q$

$B e_1 = Q^H A Q e_1 = Q^H A x_1 = Q^H \lambda_1 x_1 = \lambda_1 Q^H x_1 = \lambda_1 e_1$

prima colonna di B

$\Rightarrow B = \begin{bmatrix} \lambda_1 & * \\ 0 & A_2 \end{bmatrix}$



Abbiamo:

$$A = Q B Q^H = Q \begin{bmatrix} \lambda_1 & \underline{c}^H \\ \underline{0} & U_1 A_2 U_1^H \end{bmatrix} Q^H$$

$\downarrow$ 
← è ancora unitaria

$$= Q \begin{bmatrix} 1 & \underline{0}^H \\ \underline{0} & U_1 \end{bmatrix} \begin{bmatrix} \lambda_1 & \underline{c}^H U_1 \\ \underline{0} & A_2 \end{bmatrix} \begin{bmatrix} 1 & \underline{0}^H \\ \underline{0} & U_1^H \end{bmatrix} Q^H$$

prodotto di unitarie = U
= U^H

$\Rightarrow \exists U$  unitaria t.c.  $A = U \Lambda U^H$   $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$   
e colonne di  $U$  base autovettori ortonormali

dim:  $A = U T U^H$  per il teo sopra

$$T = U^H A U \rightarrow T^H = U^H A^H U = U^H A U = T$$

( $A^H = A$  perché  $A$  è Hermit.)

Ma  $T = T^H$  +  $T$  triangolare  $\Rightarrow T$  è diagonale e gli autovalori sono reali

→ posso scrivere  $A$  come nella Tesi e segue subito che le colonne di  $U$  sono vettori di autovettori (ortonormali perché unitaria).

## 1.4 RAGGIO SPETTRALE

def:  $S(A) = \max_i |d_i(A)|$

OSS: il raggio spettrale NON è una norma

$$\text{dim: } A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} \quad A+B = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$S(A) = 0 \quad S(B) = 0 \quad S(A+B) = 1$$

$$S(A+B) > S(A) + S(B)$$

TEO:  $\forall \|\cdot\|$  norma matriciale vale  $\rho(A) \leq \|A\|$

dim, Sia  $\lambda$  autoval  $A$  che  $\bar{\lambda}$  t.c.  $|\lambda| = S(A)$  e v autovett corrisp

$$V = \begin{bmatrix} \underline{v} & 1 & - & 1 & \underline{v} \end{bmatrix}$$

Quindi  $AV = \lambda V$ . (\*)

Abbiamo:

per proprietà submultiplic.

$$|\lambda| \|V\| = \|\lambda V\| = \|AV\| \leq \|A\| \cdot \|V\| \quad \Rightarrow \quad |\lambda| \leq \|A\| \quad \text{perché } \|V\| \neq 0$$

TEO:  $A \in \mathbb{C}^{n \times n} \Rightarrow \forall \epsilon > 0 \exists \|\cdot\|$  n.m. indotta t.c.

$$S(A) \leq \|A\| \leq S(A) + \varepsilon$$

d'um; Considero la forma di Jordan  $\rightarrow A = T J T^{-1}$

J diag. a blocchi, con blocchi del tipo  $C_i^{(s)} = \begin{bmatrix} \lambda_i & 1 \\ & \ddots \\ & & \lambda_i \end{bmatrix}$



Definiamo  $E := \text{diag}(1, \varepsilon, \dots, \varepsilon^{n-1})$

$$E^{-1} J E = E^{-1} (T^{-1} A T) E \quad \text{è diagonale a blocchi} \\ (\text{perché } E \text{ è diag. diag. blocchi. diag})$$

Come sono fatti i blocchi?

Diagram illustrating the structure of the matrix  $E$  and its effect on the matrix  $A$ . The matrix  $E$  is a block diagonal matrix with blocks of size  $\varepsilon^i$ . The matrix  $A$  is transformed into a block diagonal matrix  $D_i^{(J)} = \begin{bmatrix} \lambda_i & \varepsilon & & \\ & & \varepsilon & \\ & & & \varepsilon \\ & & & & \lambda_i \end{bmatrix} = C_i^{(J)}$ .

$$\|E^{-1} J E\|_\infty = \max_i \|D_i^{(J)}\|_\infty = \max_i (|\lambda_i| + \varepsilon) \leq \beta(A) + \varepsilon$$

il max sulla matrice è uguale al max delle norme  $\infty$  dei blocchi perché quando facciamo il max su riga (colonna) le "righe non si intersecano"

$$\text{Ma } \|E^{-1} J E\|_1 = \|E^{-1} (T^{-1} A T) E\|_1 = \|A\|_1 \quad \text{norma matriciale indotta}$$

PROP:  $\|\cdot\|_* : \mathbb{C}^{n \times n} \rightarrow \mathbb{R}^+ \cup \{0\}$  t.c.  $\|A\|_* = \|S^{-1} A S\|_\infty$  è norma matriciale indotta.

dim: (traccia)

- $\|x\|_* = \|S^{-1} x\|_\infty$  è norma vettoriale
- si considera la n.m. indotta da  $\|\cdot\|_*$  e si verifica che è quella dell'enunciato.

CONSEGUENZA:  $\beta(A) = \inf_{\| \cdot \| \text{ n.m. i. }} \|A\|$  ma in generale non è un min

OSS: Se gli autoval. che realizzano  $\beta(A)$  hanno  $\mu_3 = \mu_4$   
 $\Rightarrow \exists \|\cdot\| \text{ n.m. indotta t.c. } \beta(A) = \|A\|_*$

"dim: Qui (\*\*\*) per le roggia spett. (autoval max) posso non mettere  $\pm \varepsilon$  perché  $D_i^{(J)}$  è diag dato che  $\mu_3 = \mu_4$

Questa condizione vale sicuramente se  $A$  è diagonalizzabile. Inoltre se  $A$  è Hermitiana so che  $\|\cdot\|_* = \|\cdot\|_2$

## 1.5 SISTEMI LINEARI

PROBLEMA:  $Ax = b$   $A \in \mathbb{C}^{n \times n}$  non singolare,  $x, b \in \mathbb{C}^n$

TEO:  $Ax = b$   $A \in \mathbb{C}^{m \times n}$   $b \in \mathbb{C}^m$   $x \in \mathbb{C}^n$   
 consistente se  $\text{rang}(A/b) = \text{rang}(A)$

Possiamo formulare il nostro problema così:

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ , valore  $f$  in  $x$  assegnato

Abbiamo 3 tipi di errore:

- ERRORE INERENTE  $\rightarrow$  condizionamento  $P_0$
- ERRORE ALGORITMICO  $\rightarrow$  stabilità metodo / algoritmo
- ERRORE ANALITICO  $\rightarrow$  convergenza



L'errore inerente è generato dagli errori di rappresentazione dei dati. Devo vedere se questi errori di rappresentazione vengono amplificati! (se sì MAL condizionato)  $x \mapsto f(x)$

L'errore algoritmico è generato da errori nelle operazioni in aritmetico floating point  $f(x) + f(y) \neq f(x+y)$  se no devo forse il float  
(N.B visto che l'algoritmo lo posso cambiare devo usare metodi stabili)

L'errore analitico ce l'ho ad esempio quando passo da un problema continuo a uno discreto, o quando invece di fare un limite faccio un troncamento.

Quindi esaminiamo bene il nostro problema:

$$\begin{array}{l} x \mapsto f(x) \\ f \mapsto g \end{array} \quad \begin{array}{l} \text{(IN)} \\ \text{(AN)} \end{array} \quad \begin{array}{l} \text{f è approssimata se focessi;} \\ \text{bene i conti} \end{array} \quad \left\{ \begin{array}{l} f(x) \mapsto \psi(f(x)) \\ \downarrow \text{valore da calcolare} \end{array} \right. \quad \begin{array}{l} \downarrow \text{valore calcolato} \end{array}$$

$$\Rightarrow E_{\text{tot}} = \frac{\psi(f(x)) - f(x)}{f(x)} \stackrel{(\text{teo})}{=} E_{\text{IN}} + E_{\text{ALG}} + E_{\text{AN}}$$

Quindi posso considerare i 3 errori separatamente per avere l'errore totale.

$$E_{\text{AN}} = \frac{g(x) - f(x)}{f(x)}$$

$$E_{\text{IN}} = \frac{f(f(x)) - f(x)}{f(x)}$$

$$E_{\text{ALG}} = \frac{\psi(f(x)) - g(f(x))}{g(f(x))} \quad \leftarrow \text{qui vede solo } g \text{ e } f(x), \text{ non sa nemmeno che c'è una } f \text{ e una } x$$

\* Se  $f$  è sufficientemente regolare

$$E_{\text{IN}} = \sum_{i=1}^n c_i E_{x_i}$$

$$\text{con } E_{x_i} = \frac{f(x_i) - x_i}{x_i}, \quad |E_{x_i}| \leq \forall i \quad \text{e} \quad c_i = \frac{x_i}{f(x)} \cdot \frac{\partial f}{\partial x_i} \quad \text{COEFFICIENTI DI AMPLIFICAZIONE}$$

Sia ora  $G: \mathbb{R}^n \rightarrow \mathbb{R}^m$

$$G(\tilde{b}) = G(b) + J_G(b) \delta b \quad \text{con } \tilde{b} = b + \delta b$$

Abbiamo:

$$\|E_r\| = \frac{\|G(\tilde{b}) - G(b)\|}{\|G(b)\|} \leq \frac{\|J_G(b)\| \|\delta b\|}{\|G(b)\|} = \frac{\|J_G(b)\| \|b\|}{\|G(b)\|} \cdot \frac{\|\delta b\|}{\|b\|}$$

Ricordiamo che noi vogliamo risolvere  $Ax = b$  quindi  $x = A^{-1}b$  quindi nel nostro caso

$$G(b) = A^{-1}b$$

Quindi:

$$\frac{\|J_G(b)\| \|b\|}{\|G(b)\|} = \frac{\|A^{-T}\| \|b\|}{\|A^{-1}b\|} = \frac{\|A^{-T}\| \|Ax\|}{\|x\|} \leq \frac{\|A^{-T}\| \|A\| \|x\|}{\|x\|} = K(A) \quad \text{NUMERO DI CONDIZIONAM.}$$

In conclusione, chiamato  $E_x = \frac{\|\delta x\|}{\|x\|}$ , abbiamo ottenuto che

$$E_x \leq K(A) \cdot E_b$$

quindi  $K(A)$  rappresenta di quanto viene amplificato l'errore.



def: il **NUMERO DI CONDIZIONAMENTO** è  $K(A) = \|A\| \cdot \|A^{-1}\|$   
 Per convenzione  $K(A) = +\infty$  se  $A$  è singolare.

OSS:  $K(A) \geq 1$

Infatti:  $K(A) = \|A\| \cdot \|A^{-1}\| \geq \|AA^{-1}\| = \|I\| \geq 1$  (\*)

(\*) perché qualsiasi sia la norma usata:

$$\|I\| = \|I \cdot I\| \leq \|I\| \cdot \|I\| \Rightarrow 1 \leq \|I\|$$

def: **CONDIZIONAMENTO SPETTRALE**:  $K_2(A) = \|A\|_2 \cdot \|A^{-1}\|_2$

Ricordiamo che abbiamo visto che

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^H A)} = \sqrt{\lambda_{\max}(A^H A)}$$

La matrice  $A^H A$  è <sup>(hermitiana)</sup> simmetrica, semidef. positiva, infatti:

$$\forall x \neq 0 \quad x^H A^H A x = y^H y = \|y\|_2^2 \geq 0$$

( $x \neq 0$  e  $A$  non singolare  $\Rightarrow y \neq 0$  dove  $y = Ax$  dove  $y$  è vettore def. pos.)

$$\begin{aligned} \|A^{-1}\|_2 &= \sqrt{\lambda_{\max}(A^{-H} A^{-1})} = \sqrt{\lambda_{\max}((A^H A)^{-1})} \\ &= \sqrt{\max_i \left( \frac{1}{\lambda_i(A^H A)} \right)} = \sqrt{\frac{1}{\lambda_{\min}(A^H A)}} = \sqrt{\frac{1}{\lambda_{\min}(A^H A)}} \end{aligned}$$

perché  $A^H A \sim A^H A$

Quindi ora fine otteniamo:

$$K_2(A) = \sqrt{\frac{\lambda_{\max}(A^H A)}{\lambda_{\min}(A^H A)}} \quad \text{se } A \text{ è non singolare}$$

Se inoltre  $A$  è hermitiana  $\|A\|_2 = \rho(A)$  e quindi:

$$K_2(A) = \frac{\max_{i=1, \dots, n} |\lambda_i(A)|}{\min_{i=1, \dots, n} |\lambda_i(A)|} \quad \lambda_i(A) \in \mathbb{R} \quad (\text{perché } A \text{ hermit.})$$

Se  $A$  è anche hermitiana definita positiva

$$K_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)} \quad \lambda_i(A) \in \mathbb{R}^+$$

## 1.6 STIMA di $K(A)$

$$\forall \|\cdot\| \text{ n.m. (indotta)} \quad \rho(A) \leq \|A\|$$

$$\Rightarrow K(A) = \|A\| \cdot \|A^{-1}\| \geq \rho(A) \cdot \rho(A^{-1}) = \frac{\max |\lambda_i(A)|}{\min |\lambda_i(A)|}$$

Quindi anche usando una norma generica abbiamo una stima di  $K(A)$  che coinvolge solo gli autovalori di  $A$ .

Questa stima è dal basso, quindi ci dice quanto è almeno mal condizionato il problema.



## 1.7 METODI di FATTORIZZAZIONE

Il nostro problema è sempre risolvere  $Ax = b$  con  $A \in \mathbb{Q}^{n \times n}$ ,  $x, b \in \mathbb{C}^n$ ,  $A$  non singolare.

Abbiamo dei casi banali:

- $A = D$  diagonale  $\rightarrow x_i = b_i / a_{ii}$
- $A = L$  triang. inf  $\rightarrow$  ris. forward
- $A = U$   $\rightarrow$  ris. backward  $x_i = \frac{(b_i - \sum_{j=i+1}^n a_{ij} x_j)}{a_{ii}}$

Metodi diretti  $\rightarrow$  fattorizzazione di  $A$   
 $A = BC$

$$\Rightarrow BCx = b$$

pongo  $y = Cx \rightarrow$  Risolvo (I)  $By = b$   
 (II)  $Cx = y$

Chiaramente  $B$  e  $C$  devono essere matrici comode.

Abbiamo vari metodi di fattorizzazione:

- $LU$  con  $\text{diag}(L) = I$  Gauss  $\frac{n^3}{3}$
  - $LL^H = R^H R$  con  $\text{diag}(L^R) > 0$  Cholesky [solo se  $A$  è hermit. def. positiva]  $\frac{n^3}{6}$
- (Questi due metodi sono praticamente la stessa cosa se  $A$  è definita positiva, quale il costo comp. è lo stesso in quel caso.)
- $QR$  con  $Q$  unitaria e  $R$  triang. sup.  $2 \frac{n^3}{3}$

**OSS:** Risolvere  $Qy = b$  è comodo lo stesso perché equivale a  
 $y = Q^H b$  dato che è unitaria

(Il costo è il doppio di Gauss dato che invece di risolvere un sistema triangolare devo fare un prodotto matrice vettore)  $\leftarrow$  circa

**OSS:** se  $A$  è reale allora anche le matrici delle fattorizzazioni lo sono

**OSS:**  $LU$  non è detto che lo possa sempre fare, ma se c'è è unica.  
 $QR$  c'è sempre ma non è detto sia unica, ma questa cosa può fare anche comodo.

## 1.8 FATTORIZZAZIONE LU

A matrice da fattorizzare  $\rightarrow A^{(n)} = A$

Pongo  $A^{(k+1)} = M_k A^{(k)}$  con  $M_k$  matrice elementare di Gauss al passo  $k$   $k = 1, \dots, n-1$

$$M_k = \left[ \begin{array}{ccc|ccc} I_{k-1} & & & & & \\ & 1 & & & & \\ & 0 & -m_{kk} & & & \\ & & & 1 & & \\ & & & & \ddots & \\ & & & & & 1 \end{array} \right] \quad \text{con} \quad m_{ik} = \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}}$$

Ma possiamo scrivere  $M_k$  anche nella forma:

$$M_k = I - \underline{m}_k \cdot \underline{e}_k^T \quad \text{con} \quad \underline{m}_k = \begin{bmatrix} 0 \\ \vdots \\ m_{ik} \\ \vdots \end{bmatrix}_{i=k}^T$$

$\uparrow$  DIAG  $\quad \cdot \quad \underline{e}_k^T = \text{matrice } 1$

Abbiamo:

$$(M_{n-1} \dots M_1) A = U$$

$$(M_{n-1} \dots M_1)^{-1} = M_1^{-1} \dots M_{n-1}^{-1} = L$$



**TEO:**  $A \in \mathbb{C}^{n \times n}$ . Se  $A_k$  sottomatrice principale di testa  $k$  non singolare  $\forall k = 1, \dots, n \quad \exists! A = LU$

**TEO:**  $A \in \mathbb{C}^{n \times n}$ .  $\exists P$  matrice di permutazione t.c.  $PA = LU$

**N.B.:** Nel metodo in pratica si usa sempre la permutazione perché così viene meglio per i calcoli.  
(Pivot parziale, Pivot totale)  
→ **EFFETTO STABILIZZANTE:**

Pivot parziale:



$$1. |M_{ik}| \leq 1 \quad (\text{dato che } a_{kk} \leq a_{kk})$$

$$2. a_m^{(n)} \leq 2^{n-1} a_m^{(1)}$$

Pivot totale in pratica non si usa quasi mai, perché è meglio ma non serve di solito.

**ESEMPIO:**  $A = \begin{bmatrix} \varepsilon & 1 \\ 1 & 0 \end{bmatrix} \quad b = \begin{bmatrix} 1+\varepsilon \\ 1 \end{bmatrix} \quad \varepsilon > 0 \quad x_{esatta} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

Il problema è ben condizionato:  $A^{-1} = \begin{bmatrix} 0 & 1 \\ 1 & -\varepsilon \end{bmatrix}$

$$K_\infty(A) = \|A\|_\infty \|A^{-1}\|_\infty = (1+\varepsilon)^2 \approx 1 \quad \text{per } \varepsilon \text{ piccolo}$$

Pero' se faccio LU ottengo:

$$A = LU = \begin{bmatrix} 1 & 0 \\ \frac{1}{\varepsilon} & 1 \end{bmatrix} \begin{bmatrix} \varepsilon & 1 \\ 0 & -\frac{1}{\varepsilon} \end{bmatrix} \quad \text{Per } \varepsilon \text{ piccolo esplodono!!}$$

→ bisogna fare il pivot parziale!

Ad esempio si può verificare che in  $A(10^{-3})$  con  $\varepsilon = 0.3 \cdot 10^{-3}$  si ottiene  $\bar{x} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$  (cifre di mantissa)

## 1.9 FATTORIZZAZIONE QR

**def:** **MATRICE ELEMENTARE di HOUSEHOLDER:**

$$P = I - \beta v \cdot v^H \quad \beta \in \mathbb{R}, v \in \mathbb{C}^n \quad v \neq 0 \rightarrow v v^H \text{ è matrice}$$

**PROP:**  $P$  è hermitiana.

$P$  è unitaria se  $\beta = \frac{2}{\|v\|_2^2}$

(i.e. dato un vettore  $v$  posso associargli un'unica matrice unitaria e hermit. in questa forma)

dim:  $P^H = I - \beta (v^H)^H \cdot v^H = I - \beta v \cdot v^H = P$

Deve valere  $I = P P^H = (I - \beta v v^H)(I - \beta v v^H)$

$$\begin{aligned} & \stackrel{P=P^H}{=} I - 2\beta v v^H + \beta^2 v v^H v v^H \\ & \stackrel{\|v\|_2^2}{=} I - 2\beta v v^H + \beta^2 \|v\|_2^2 v v^H \end{aligned}$$

Questo vale  $\Leftrightarrow -2\beta v v^H + \beta^2 \|v\|_2^2 v v^H = 0$

$$(-2 + \beta \|v\|_2^2) \beta v v^H = 0$$

Se  $\beta \neq 0$  (caso interessante) deve essere  $-2 + \beta \|v\|_2^2 = 0$



PROP: Sia  $P = P(\underline{v}) = I - \beta \underline{v} \underline{v}^H$  con  $\beta = \frac{2}{\|\underline{v}\|_2^2}$ . Allora:

•  $\forall \underline{x} \in \underline{v}^\perp = \{\underline{x} : \underline{x}^H \underline{v} = 0\}$  vale  $P\underline{x} = \underline{x}$  (i.e.  $P$  è identità su  $\underline{v}^\perp$ )

•  $P(\underline{v}) \underline{v} = -\underline{v}$  (i.e. è una riflessione nel stsp. gen da  $\underline{v}$ )

$$\text{dim: } P\underline{x} = (I - \beta \underline{v} \underline{v}^H) \underline{x} = \underline{x} - \underbrace{\beta \underline{v} \underline{v}^H \underline{x}}_{=0} = \underline{x}$$

$$P(\underline{v}) \underline{v} = (I - \beta \underline{v} \underline{v}^H) \underline{v} = \underline{v} - \beta \|\underline{v}\|_2^2 \cdot \underline{v} \\ = \underline{v} - 2\underline{v} = -\underline{v}$$

$$\text{ma } \beta = \frac{2}{\|\underline{v}\|_2^2}$$

TEO:  $\forall \underline{x} \in \mathbb{C}^n \setminus \{0\} \exists P \in \mathbb{C}^{n \times n}$  m. e. di Householder t.c.

24/10/2012

$$P\underline{x} = \alpha \underline{e}_1 \quad \alpha \in \mathbb{C} \text{ opportuno, } \underline{e}_1 = [1, 0, \dots, 0]^T$$

dim: (I) troviamo  $\alpha \in \mathbb{C}$

(II) troviamo  $\underline{v} \in \mathbb{C}^n$  che individua  $P$

(I) Sia  $\alpha = |\alpha| e^{i\theta}$  (lo scrivo in forma exp)

Vogliamo che  $P\underline{x} = \alpha \underline{e}_1$  e quindi deve essere

$$\|P\underline{x}\|_2 = \|\alpha \underline{e}_1\|_2 = |\alpha|$$

Ma  $P$  è unitaria e quindi  $\|P\underline{x}\|_2 = \|\underline{x}\|_2 \Rightarrow |\alpha| = \|\underline{x}\|_2$  (\*)

Deve anche essere  $\underline{x}^H P\underline{x} = \underline{x}^H \alpha \underline{e}_1$  (analogamente visto che  $P\underline{x} = \alpha \underline{e}_1$ )

Abbiamo  $\underline{x}^H \alpha \underline{e}_1 = \alpha \bar{x}_1 \in \mathbb{R}$  (\*\*) perché  $P$  è hermitiana

$$\underline{x}^H \underline{e}_1$$

$$\underline{y}^H P \underline{y} = (\underline{y}^H P \underline{y})^H = \underline{y}^H P^H \underline{y} = \underline{y}^H P \underline{y} \in \mathbb{R}$$

$\downarrow$   
 $P = P^H$

$$x_1 = |x_1| e^{i\varphi} \neq 0 (**)$$

Abbiamo:

$$\alpha \bar{x}_1 = |\alpha| e^{i\theta} |x_1| e^{-i\varphi} = \|\underline{x}\|_2 |x_1| e^{i(\theta-\varphi)} \in \mathbb{R}$$

$$\Rightarrow \theta - \varphi \in k\pi \quad \text{e quindi} \quad \theta = \varphi + \delta\pi \quad \delta \in \{0, 1\}$$

e quindi:

$$\alpha = |\alpha| e^{i\theta} = \|\underline{x}\|_2 e^{i(\varphi + \delta\pi)} = \pm \|\underline{x}\|_2 e^{i\varphi} = \pm \|\underline{x}\|_2 \frac{x_1}{|x_1|}$$

nella pratica userei  
il segno -

Se  $x_1 = 0$  per convenzione si pone  $\alpha = -\|\underline{x}\|_2$  (ho libertà di scelta per  $\theta$ )  
in modo che sia soddisfatta (\*).

è un prodotto scalare quindi è un numero e posso spostarlo davanti

$$(II) P\underline{x} = (I - \beta \underline{v} \underline{v}^H) \underline{x} = \underline{x} - \underbrace{\frac{2}{\|\underline{v}\|_2^2} (\underline{v}^H \underline{x}) \underline{v}}_{\hat{\beta} \underline{v}} = \underline{x} - \hat{\beta} \underline{v}$$

$$\text{di vuole } \hat{\beta} \underline{v} = \underline{x} - P\underline{x} = \underline{x} - \alpha \underline{e}_1$$

Questa relazione non è tanto bella perché  $\hat{\beta}$  dipende da  $\underline{v}$  e quindi la relazione non è lineare! (...)

$$\text{OSS: } P[\underline{x}\underline{v}] := I - \frac{2}{\|\underline{x}\underline{v}\|_2^2} (\underline{v}\underline{v}^H) (\underline{x}\underline{x}^H) = I - \frac{2}{\|\underline{x}\|^2 \|\underline{v}\|_2^2} \|\underline{x}\|^2 \underline{v} \underline{v}^H = P[\underline{v}]$$

quindi la matrice generata da  $\underline{v}$  e da  $\underline{x}\underline{v}$  è la stessa!!

(...) Quindi noi possiamo prendere  $\hat{\beta} = 1$  tutto lo scalatura viene mangiata (basta che sia solo  $\hat{\beta} \neq 0$ )

$$\Rightarrow \text{prendo } \underline{v} = \underline{x} - \alpha \underline{e}_1$$



Ora il punto è: posso davvero escludere che  $\hat{\beta} = 0$ ?

Per def  $\hat{\beta} = 0$  sse  $\underline{y}^H \underline{x} = 0$  sse  $\underline{x} \in \underline{V}^\perp$  (\*)

Dato che  $\underline{x}$  ce l'ho e  $\underline{V}$  lo sto scegliendo basta che lo scelgo in modo che (\*) non valga!

Osserviamo che la nostra scelta va bene:

$$\begin{aligned}\underline{V}^H \underline{x} &= (\underline{x} - \alpha \underline{e}_1)^H \underline{x} = \underline{x}^H \underline{x} - \alpha \underline{e}_1^H \underline{x} \\ &= \|\underline{x}\|_2^2 - \alpha x_1 = \|\underline{x}\|_2^2 + \|\underline{x}\|_2^2 \frac{x_1}{|x_1|} \cdot x_1 \quad \rightarrow \bar{x}_1 x_1 = |x_1|^2 \\ &= \|\underline{x}\|_2^2 + |x_1| \|\underline{x}\|_2^2 \neq 0 \quad (> 0)\end{aligned}$$

Ricapitolando:

$$\alpha = \begin{cases} \pm \|\underline{x}\|_2 \frac{x_1}{|x_1|} & \text{se } x_1 \neq 0 \\ -\|\underline{x}\|_2 & \text{se } x_1 = 0 \end{cases}$$

$$\underline{V} = \underline{x} - \alpha \underline{e}_1 = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

Così  $v_1$ ? Prendiamo  $\alpha = -\|\underline{x}\|_2 \frac{x_1}{|x_1|}$

$$v_1 = x_1 - \alpha = x_1 - \left(-\|\underline{x}\|_2 \frac{x_1}{|x_1|}\right) = x_1 + \frac{x_1}{|x_1|} \|\underline{x}\|_2 = \frac{x_1}{|x_1|} (|x_1| + \|\underline{x}\|_2)$$

quindi scelgo  $\alpha$  con segno "-" per questo  $\leftarrow$  NON RISCHIO DI AVERE ERRORI DI CANCELLAZIONE

OSS: Se  $\underline{x} \in \mathbb{R}^n \rightarrow \alpha \in \mathbb{R}$  e  $\underline{V} \in \mathbb{R}^n$  per costruzione  
 $\Rightarrow P$  è ortogonale e  $P \in \mathbb{R}^{n \times n}$

Ricordiamo che vogliamo scrivere  $A = QR$ . Allora cosa facciamo in generale al passo  $k$ -esimo. (Abbiamo visto ora il primo)

Passo  $k$ -esimo:  $\underline{Q}_k$  sottocolumna  $k$ -esima  $\in \mathbb{C}^{n \times 1}$

$$\underline{V}_k \in \mathbb{C}^n, \underline{V}_k = \begin{bmatrix} 0 \\ \vdots \\ v_{k,k} \\ \vdots \\ v_{n,k} \end{bmatrix} \quad |k-1$$

Quindi: procedo facendo  $P^{(n-1)} \dots P^{(2)} P^{(1)} A = \begin{bmatrix} \text{---} & \text{---} & \text{---} \\ 0 & \text{---} & \text{---} \\ \vdots & \ddots & \vdots \\ 0 & \text{---} & \text{---} \end{bmatrix} = R$

Ma allora abbiamo:

$$\begin{aligned}A &= (P^{(n-1)} \dots P^{(1)})^{-1} R \\ &= P^{(1)-1} \dots P^{(n-1)-1} R = P^{(1)} \dots P^{(n-1)} R = QR \\ &\quad \text{perché } P^{-1} = P^H = P \quad \uparrow \text{ prodotto di matrici unitarie e unitarie}\end{aligned}$$

Vediamo ora la NON UNICITÀ della scomposizione  $QR$ .

$$S = \text{diag}(\theta_1, \dots, \theta_n) \quad \text{unitaria} \quad |\theta_i| = 1 \quad \forall i \quad \text{MATRICE DI FASE}$$

$$A = QR = QSS^H R = (QS)(S^H R) = \tilde{Q} \tilde{R}$$

$\tilde{Q}$  è unitaria perché prodotto di unitarie.

$\tilde{R}$  è triangolare sup perché  $S^H$  è diagonale.



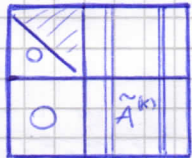
Quindi ho infiniti modi per fattorizzare col metodo QR.

Abbiamo che il costo del metodo QR è  $O(n^3)$  ed è doppio rispetto a quello LU (è sempre  $O(n^3)$  ma doppio)

Dato che il costo è doppio di solito si usa il metodo LU. La fattoriz. QR si usa però per esempio per il calcolo del rango di una matrice o all'interno del metodo SVD...

Esiste una tecnica di pivot anche per il metodo QR. A che cosa serve? (Visto che TUTTA la fattorizzazione esiste) Generalmente serve per avere un ordinamento negli elementi della diagonale della fatt.

MASSIMO PIVOT per COLONNE in QR:

Posso  $k$   $A^{(k)} =$    $k-1$  cerco la colonna di  $A^{(k)}$  con norma 2 max e scambio le colonne. ( $\pi^{(k)}$  matrice di scambio)

Ala fine ho  $P^{(n-1)} \dots P^{(1)} A \underbrace{\pi^{(1)} \dots \pi^{(n-1)}}_{=\pi} = R$    
 multiplico di qui perché il lavoro sulle colonne

$$\Rightarrow A\pi = QR$$

(\*) In questo modo la matrice  $R$  è t.c.  $|r_{11}| \geq |r_{22}| \geq \dots \geq |r_{nn}|$    
 dim; Nella notazione del TEO di esistenza di  $P^{(i)}$  abbiamo

$r_{ii} = q_i$    
 Avevamo visto che  $|q_1| = \|a_1\|_2$  con  $A\pi^{(1)} = [a_1 | a_2 | \dots | a_n]$    
 per costruzione di  $P^{(1)}$

E' quindi evidente che  $|q_1| = \|a_1\|_2 \geq \|a_i\|_2 \quad i=2, \dots, n$  perché ho scelto  $\pi^{(1)}$    
 So anche però che  $\|q_i\|_2 = \|P^{(i-1)} a_i\|_2$  dato che  $P^{(i-1)}$  è unitaria.

$$\text{Chiamiamo } A^{(2)} = P^{(1)} A \pi^{(1)} \pi^{(2)}$$

Ovviamente  $\|P^{(1)} a_i\|_2 \geq \|\tilde{a}_i\|_2$  perché  $\tilde{a}_i$  è uguale a  $P^{(1)} a_i$  ma con delle componenti in meno



Quindi alla fine arrivo che  $|q_1| \geq \|\tilde{a}_2\|_2 = |q_2|$

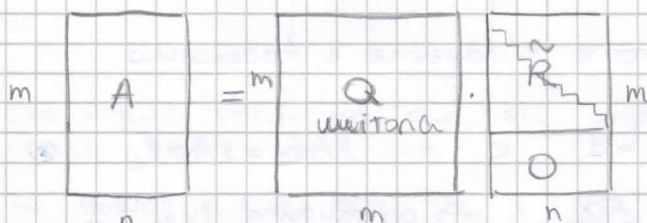
poi si itera per gli altri  $\tilde{a}_i$ . Dato che  $q_i = r_{ii}$  ottengo la tesi.

QR per calcolo del RANGO   
 Vediamo come usiamo questo metodo per calcolare il rango.

$A \in \mathbb{C}^{m \times n}$ . Se  $\text{rk}(A) = p < n \Rightarrow$  dopo  $p$  passi si ferma   
 applico QR con pivot max  $A\pi = QR$    
  $\det(R_1) \neq 0$   $R = \begin{bmatrix} R_1 & S \\ 0 & 0 \end{bmatrix} \begin{matrix} p \\ n-p \end{matrix}$

FATTORIZZAZIONE QR PER MATRICI RETTANGOLARI

$$A \in \mathbb{C}^{m \times n} \quad m > n \quad A = QR$$



$\tilde{R}$  triang. superiore  $\in \mathbb{C}^{n \times n}$

OSS: il procedimento è lo stesso cambia solo che ho  $n$  passi perché devo azzerare anche l'ultima colonna. (Se è quadrata non ho niente da azzerare nell'ultima colonna)

$$\text{COSTO} = n^2 \left( m - \frac{n}{3} \right)$$